

### **Amendments to the Claims**

**This listing of claims will replace all prior versions, and listings, of the claims:**

1. (previously presented) A server system comprising:
  - at least two scaleable tiers of server machines;
  - a server pool including plural spare server machines;
  - means for computing an average response time for the server system to respond to transaction requests at the two scaleable tiers of server machines; and
  - means for increasing a number of server machines processing transactions for each of the two scaleable tiers of server machines by allocating the spare server machines to process a portion of the transactions, wherein the spare server machines are allocated to process a portion of the transactions when the average response time for the server system to respond to the transaction requests is greater than or equal to a specified average response time.
2. (previously presented) The server system of claim 1 further comprising means for determining costs associated with allocating the number of server machines at each of the two scaleable tiers of server machines.
3. (previously presented) The server system of claim 2 wherein said means for determining further comprises means for minimizing costs associated with allocating an optimized number of server machines at each of the two scaleable tiers of server machines.
4. (previously presented) The server system of claim 3 wherein said means for minimizing comprises:
  - means operatively coupled to said server system for receiving input parameters and for solving:

$$\sqrt{\gamma} = \frac{\sum_{i=1}^n \sqrt{h_i s_i u_i}}{T - \sum_{i=1}^n s_i};$$

where:  $\gamma$  is a shadow price of the average response time;  $h_1, h_2, \dots, h_n$  are weights reflecting a cost of different types of servers located at each of the two scaleable tiers of server machines;  $s$  is an average service time;  $u$  is a measured average utilization rate expressed in a single-machine percentage; and  $T$  is the average response time.

5. (previously presented) The server system of claim 1, wherein the average response time is determined by examining a time that the transaction requests are pending at each of the two scaleable tiers of server machines.

6. (previously presented) The server system of claim 1 further comprising:

a contractual relationship between a system operator and at least one contracting party; and

means for adjusting prices charged by said system operator to at least one third party in response to a change in an allocation of server machines in each of the two scaleable tiers of server machines.

7. (previously presented) The server system of claim 1 wherein said means for computing further comprises a non-iterative queuing model for predicting the average response time for the server system in response to measured arrival rates of transaction requests into each of the two scaleable tiers of server machines, an average service demand at each of the two scaleable tiers of server machines, and a number of servers allocated to each of the two scaleable tiers of server machines.

8. (previously presented) A method for allocating a server machine to at least two tiers of a server system, said method comprising:

computing an expected average response time as a function of transaction requests and an amount of resources allocated to each of the two tiers of the server system;

determining whether an optimization problem is feasible;

computing a lower bound and an upper bound on a number of server machines at each of the two tiers of said server system required to meet the average response time;

computing a solution specifying a number of server machines allocated to each of the two tiers of said server system;

computing an average time that transaction requests are pending at each of the two tiers;

automatically increasing the number of server machines allocated to one of the two tiers at a point in time when the average time the transaction requests are pending at the one of the two tiers is greater than or equal to a pre-determined limit.

9. (previously presented) The method of claim 8 wherein said computing an expected average response time further comprises:

obtaining at least one input value for an average arrival rate of transaction requests into each of the two tiers of said server system;

obtaining at least one input value for an average service demand at each of the two tiers of said server system; and

obtaining at least one input value for the number of server machines allocated at each of the two tiers of said server system.

10. (canceled)

11. (previously presented) An assembly for allocating server machines in a server system comprising:

at least two tiers of server machines;

a pool of spare server machines that process transactions for the two tiers of server machines;

means for computing an average response time for said two tiers of server machines to respond to a plurality of transaction requests; and

means for increasing and decreasing a number of server machines from said pool that process transactions for said two tiers of server machines when average response times for processing transactions at the two tiers of server machines exceed a specified average response time.

12. (previously presented) The assembly of claim 11, wherein the average response time is determined by examining a time that the transaction requests are pending at the two tiers of server machines.

13. (previously presented) The assembly of claim 11 further comprising:

a contractual relationship between a system operator and at least one contracting party; and

means for adjusting prices charged by said system operator to said at least one contracting party in response to a change in an allocation of server machines in said two tiers of server machines.

14. (previously presented) The assembly of claim 11 wherein said means for computing further comprises a non-iterative queuing model for predicting an average server system response time in response to measured arrival rates of transaction requests into said two tiers of server machines, an average service demand at said two tiers of server machines; and a number of servers allocated to said two tiers of server machines.

15. (previously presented) A server system comprising:

an open queuing network of multiple server machines with each server machine having a processor-sharing queue with a single critical resource;

at least two tiers of server machines; and

a computer-readable medium comprising instructions for:

(i) predicting an average system response time of said multiple server machines based on an arrival rate of transaction requests into each of the two tiers of server

machines averaged over all transaction request types and a number of server machines allocated at each of the two tiers of server machines;

(ii) solving a mathematical representation of an optimization objective and constraints of said server system;

(iii) determining a number of server machines for each of the two tiers of server machines in response to said predicted the average system response time; and

(iv) automatically increasing the number of server machines processing transactions for each of the two tiers of server machines at a point in time when an average time that transactions requests are pending at the two tiers of server machines exceeds a threshold.

16. (previously presented) The server system of claim 15 wherein said mathematical representation comprises:

a continuous-relaxation model of a mathematical optimization system; and  
an iterative bounding procedure.

17. (previously presented) The server system of claim 15 wherein said instructions for determining the number of server machines for each of the two tiers of server machines is in response to a predicted average system response time and at least one service level agreement (SLA) requirement.

18. (canceled)